

日本国特許庁
JAPAN PATENT OFFICE

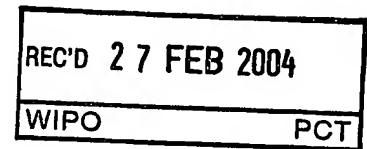
10. 2. 2004

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出願年月日
Date of Application: 2003年12月 5日

出願番号
Application Number: 特願2003-406776
[ST. 10/C]: [JP 2003-406776]



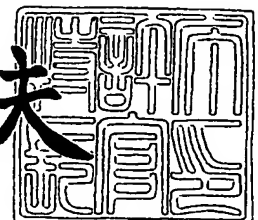
出願人
Applicant(s): 独立行政法人産業技術総合研究所

PRIORITY DOCUMENT
SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH
RULE 17.1(a) OR (b)

2004年 1月 9日

特許庁長官
Commissioner,
Japan Patent Office

今井康夫



【書類名】 特許願
【整理番号】 217-03631
【あて先】 特許庁長官殿
【国際特許分類】 G01N 33/483
G01N 33/68
【発明者】
【住所又は居所】 東京都江東区青海 2 - 4 1 - 6 独立行政法人産業技術総合研究
所臨海副都心センター内
【氏名】 富井 健太郎
【特許出願人】
【識別番号】 301021533
【氏名又は名称】 独立行政法人産業技術総合研究所
【代表者】 吉川 弘之
【電話番号】 029-861-3280
【先の出願に基づく優先権主張】
【出願番号】 特願2002-377704
【出願日】 平成14年12月26日
【提出物件の目録】
【物件名】 特許請求の範囲 1
【物件名】 明細書 1
【物件名】 図面 1
【物件名】 要約書 1

【書類名】 特許請求の範囲

【請求項 1】

タンパク質の立体構造を予測するためのタンパク質プロファイル行列間の類似性を評価するシステムであって、

前記プロファイル行列は、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントにおいて、各アミノ酸残基位置におけるアミノ酸種毎の出現確率を備えたプロファイルカラムの群から構成され、

前記類似性評価システムは、以下の手段：

- (a) 立体構造を予測したいタンパク質を含む複数のタンパク質に基づいて作成される入力プロファイル行列と、立体構造が既知である複数のタンパク質に基づいて作成される対象プロファイル行列の2つのプロファイル行列を用意する手段と、
- (b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出する手段と、
- (c) 前記相関係数からなるスコア行列を作成する手段とを含むシステム。

【請求項 2】

請求項 1 記載のシステムにより作成されたスコア行列を用いることを特徴とするタンパク質立体構造の予測システム。

【請求項 3】

コンピュータを、タンパク質の立体構造を予測するためのタンパク質プロファイル行列間の類似性を評価するシステムとして機能させるためのプログラムであって、

前記プロファイル行列は、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントにおいて、各アミノ酸残基位置におけるアミノ酸種毎の出現確率を備えたプロファイルカラムの群から構成され、

前記類似性評価システムは、以下の手段：

- (a) 立体構造を予測したいタンパク質を含む複数のタンパク質に基づいて作成される入力プロファイル行列と、立体構造が既知である複数のタンパク質に基づいて作成される対象プロファイル行列の2つのプロファイル行列を用意する手段と、
- (b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出する手段と、
- (c) 前記相関係数からなるスコア行列を作成する手段とを含むプログラム。

【請求項 4】

請求項 3 記載のプログラムを記録したコンピュータ読み取り可能な記録媒体。

【書類名】 明細書

【発明の名称】 タンパク質立体構造予測システム

【技術分野】

【0001】

本発明は、タンパク質プロファイル行列間の類似性を評価するシステムに関するものであり、より詳しくは、タンパク質の立体構造予測に好適に使用されるタンパク質プロファイル行列間の類似性の評価システムに関する。

【背景技術】

【0002】

自然界にあるタンパク質は進化の過程で選択され、特定の機能を発現するに至ったが、このタンパク質の機能はその立体構造に依存することが知られている。したがって、タンパク質の立体構造が予測できれば、その機能を予測することが可能となる。

【0003】

従来、未だ何の知見も得られていないタンパク質を調べるに際し、既に立体構造が知られているタンパク質との類似性をコンピュータによって測定することにより、タンパク質の立体構造を推論ないし予測する手法が望まれていた。このような手法の1つとして、タンパク質プロファイル行列同士を比較する方法が、有力な手法として知られている (Rychlewski L, Jaroszewski L, Li W, Godzik A. Protein Sci (2000) Feb;9(2):232-41: 非特許文献1)。

【0004】

ここで、タンパク質プロファイル行列とは、関連するタンパク質 (タンパク質ファミリーなど) におけるアミノ酸種の出現頻度を、そのアミノ酸残基位置毎に数値化して行列としたものである。この行列は、通常、以下の手順で作成される。すなわち、まず、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントを与えられ、マルチプルアライメントの各アミノ酸残基位置における20種のアミノ酸の各種類の出現数が計算される。続いて、これらの数を規格化することによって、出現確率に転換される。この時、与えられたマルチプルアライメントに含まれるメンバー内での相互のアミノ酸配列類似性に応じた重みが考慮された上で出現数が補正され、プロファイル行列が作成される。

【0005】

ここで、マルチプルアライメントとは、生物学的に相互に関連しあう複数のタンパク質のアミノ酸配列を、対応すると考えられるアミノ酸残基を揃えて並置したものをいう。マルチプルアライメントは、例えば、ある一配列を入力値として、既存のプログラムであるPSI-BLAST(Altschul et al., Nucleic Acids Res. (1997) 25(17):3389-3402: 非特許文献2)を用いて、配列データベースに検索をかけることや、生物学的に相互に関連しあう複数のタンパク質のアミノ酸配列の一群を入力値として、これも既存のプログラムであるCLUSTALW(Higgins D., Thompson J., Gibson T. Thompson J.D., Higgins D.G., Gibson T.J. (1994). Nucleic Acids Res. 22:4673-4680: 非特許文献3)を用いることで容易に作成することができる。また、立体構造比較などの結果から作成することも可能である。

【0006】

表1は、アミノ酸配列の長さ (アミノ酸残基数) が n であるタンパク質を基準として作成されたマルチプルアライメントを模式的に示したものである。なお、表1中、第1列目は個々のタンパク質の名称であり、第1行目の「1~ n 」は、マルチプルアライメントにおけるアミノ酸残基位置を示す。また、表1中のアルファベットはアミノ酸種を1文字表記したものである。

【0007】

【表1】

	1	2	3	4	5	6	7	8	...	n
20807455/14-218	M	I	D	H	T	L	L	K	...	G
19551629/13-215	I	L	D	Y	T	L	L	G	...	A
16974933/15-229	L	M	D	L	T	T	L	N	...	A
16120769/20-234	L	M	D	L	T	T	L	N	...	A

【0008】

表1の例では、例示されたアミノ酸残基位置のすべてにアミノ酸が配置されているが、アミノ酸残基位置に対応するアミノ酸残基がないとされた場合は、「・（ドット）」としてギャップを示すこともできる。表2は、表1で得られた長さがnであるマルチプルアライメントにしたがって作成されたプロファイル行列を模式的に示したものである。表2中、第1列目はアミノ酸種（ギャップを含んでいてもよい）であり、第1行目の「1～n」は、プロファイル行列におけるアミノ酸残基位置を示す。

【0009】

【表2】

AA/Pos.	1	2	3	...	n
A	0.00	0.00	0.00	...	0.71
R	0.00	0.00	0.00	...	0.00
N	0.00	0.00	0.00	...	0.00
D	0.00	0.00	0.96	...	0.00
C	0.00	0.00	0.00	...	0.00
Q	0.00	0.00	0.00	...	0.00
E	0.00	0.00	0.04	...	0.00
G	0.00	0.00	0.00	...	0.29
H	0.00	0.00	0.00	...	0.00
I	0.29	0.29	0.00	...	0.00
L	0.41	0.29	0.00	...	0.00
K	0.00	0.00	0.00	...	0.00
M	0.29	0.41	0.00	...	0.00
F	0.00	0.00	0.00	...	0.00
P	0.00	0.00	0.00	...	0.00
S	0.00	0.00	0.00	...	0.00
T	0.00	0.00	0.00	...	0.00
W	0.00	0.00	0.00	...	0.00
Y	0.00	0.00	0.00	...	0.00
V	0.01	0.01	0.00	...	0.00

【0010】

プロファイル行列中の各列は、関連する複数のタンパク質における、各アミノ酸残基位置の全アミノ酸種の確率分布を表すことになる。表3は、表2に示されたプロファイル行

列のうち、残基位置が「2」であるプロファイルカラムを模式的に示したものである。
【 0 0 1 1 】

【表 3】

2
0.00
0.00
0.00
0.00
0.00
0.00
0.00
0.00
0.00
0.00
0.29
0.29
0.00
0.41
0.00
0.00
0.00
0.00
0.00
0.00
0.01

【0012】

すなわち、表2で示されるプロファイル行列では、残基位置が2におけるアラニン (A) の補正された出現確率は0.00であり、メチオニン (M) の補正された出現確率は0.41ということになる。

【0013】

従来、2つのプロファイル行列や2つのアミノ酸配列を比較及び／又は揃えるために、ダイナミックプログラミング (Needleman SB, Wunsch CD, J Mol Biol. (1970) Mar;48(3):443-53 : 非特許文献4) が使用されてきた。アラインメントを作成する時に、比較される2つのアミノ酸配列や2つのプロファイル行列中のどの残基又はプロファイルカラムを対応付させるか (そこでは残基とギャップとの対応付も含まれる) 決定する必要があるが、その対応付のさせ方は非常に多数考えられる。ダイナミックプログラミングは、その中から類似性スコアが最大となるような対応付を自動的に効率良く見出すアルゴリズムである。そしてまた、その対応付の結果それ自体が最終的に得たいアラインメントである。

【0014】

ダイナミックプログラミングでは、通常のアミノ酸配列比較の場合は、比較される2つのアミノ酸配列、および、比較したい2つのアミノ酸配列の各々の残基ペアに対する類似性スコア (類似の度合いを示す点数) から構成されるスコア行列、プロファイル行列比較の場合は、比較される2つの代表アミノ酸配列と、比較したい2つのプロファイル行列の、各々のプロファイルカラムのペアに対する類似性スコアから構成されるスコア行列の入力を要求する。これらを入力することによって、ダイナミックプログラミングは、通常のアミノ酸配列比較の場合は、比較されるアミノ酸配列対のアラインメントとその最終スコア (類似性スコアが最大となるような最適パスを見つけることにより得られたスコア値)、プロファイル行列比較の場合は、比較される代表アミノ酸配列のアラインメント、およびその最終スコアが出力される。

【0015】

したがって、ダイナミックプログラミングを使用する手法によりプロファイル行列を比較するためには、比較したい2つのプロファイル行列の類似性を精度よく評価したスコア行列を作成する必要がある。

【0016】

2つのプロファイル間の類似の程度を示すスコア行列を算出する方法として、Rychlewskiらが開発した手法が知られている (Rychlewski et al. (2000), 9:p232-241)。これは、比較したいプロファイルカラムペア間の類似性スコアを、2つのプロファイルカラムを内積したものと定義づけて算出することにより、比較したい2つのプロファイル行列間のスコア行列を作成するものである。

【0017】

たとえば、2つのプロファイル行列、 $X = x_1 x_2 \cdots x_p \cdots x_n$ (ただし、 x_p はアミノ酸残基位置 p におけるプロファイルカラム) および $Y = y_1 y_2 \cdots y_q \cdots y_m$ (ただし、 y_q はアミノ酸残基位置 q におけるプロファイルカラム) が与えられたとき、 n 行 m 列のスコア行列の要素である、類似性スコア $D_{q,p}$ (プロファイルカラム x_p およびプロファイルカラム y_q 間の類似性スコア) は、下記の式によって与えられる。

【0018】

【数 1】

$$D_{pq} = \sum_a^j x_{pa} y_{qa}$$

[式中、 x_{pa} = プロファイルカラム x_p の要素
 y_{qa} = プロファイルカラム y_q の要素
 j = プロファイルカラムの要素数 (通常 20) である。]

【0019】

当該手法によれば、比較したい2つのプロファイルカラム間において、共にアミノ酸置換が激しくない出現残基種が非常に限られている場合には、内積した値も高い数値となるため、高い類似性スコアが与えられる事になる。このように出現残基種が非常に限られておりアミノ酸変異が激しくない高度に保存されている残基位置は、生体内での機能的あるいは、物理化学的要請から高度に保存された箇所と考えられ、生物学的にも重要な位置であると考えられている。上記手法では、このような領域はその類似性を精度良く評価することができると考えられる。

【0020】

しかしながら、上記手法では、こうした出現残基種が限られた位置を精度良く評価することができる可能性があるものの、生物学的に重要な位置であっても、モチーフ内に存在する非保存位置や、タンパク質立体構造上露出していることが重要で極性のみが重大な意義を占める位置、あるいはその逆に埋没部分に位置し疎水性のみが保存されている位置など、アミノ酸置換が激しく生起していてもその置換パターンに共通性があると考えられるような領域に関して精度良く評価することができないという問題があった。

【0021】

さらに、スコア行列の各要素 (類似性スコア) の平均値は負の値である事、標準偏差もほぼ一定値である事が望まれるため、類似性スコアに対して正規化処理を施さなければならず、煩雑であるという問題もあった。

【0022】

従って、プロファイル行列間において、保存領域のみならず、非保存領域の類似性も評価もできる、高精度かつ簡便な手法の開発が望まれていた。

【非特許文献1】 Rychlewski L, Jaroszewski L, Li W, Godzik A. Protein Sci 2000 Feb;9(2):232-41

【非特許文献2】 Altschul et al., Nucleic Acids Res. (1997) 25(17):3389-3402

【非特許文献3】 Higgins D., Thompson J., Gibson T. Thompson J.D., Higgins D.G., Gibson T.J. (1994). Nucleic Acids Res. 22:4673-4680

【非特許文献4】 Needleman SB, Wunsch CD, J Mol Biol. 1970 Mar;48(3):443-53

【発明の開示】

【発明が解決しようとする課題】

【0023】

本発明は、タンパク質の立体構造を予測するための、タンパク質プロファイル行列同士

の類似性を評価するシステムを提供することを目的とする。

【課題を解決するための手段】

【0024】

すなわち、本発明は、次のようなタンパク質プロファイル行列間の類似性評価システム、タンパク質立体構造の予測システム、コンピュータをそれらシステムとして機能させるためのプログラム、そのプログラムを記録したコンピュータ読み取り可能な記録媒体等を提供する。

【0025】

(1) タンパク質の立体構造を予測するための、タンパク質プロファイル行列間の類似性を評価するシステムであって、

前記プロファイル行列は、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントにおいて、各アミノ酸残基位置におけるアミノ酸種毎の出現確率を備えたプロファイルカラムの群から構成され、

前記類似性評価システムは、以下の手段：

(a) 立体構造を予測したいタンパク質を含む複数のタンパク質に基づいて作成される入力プロファイル行列と、立体構造が既知である複数のタンパク質に基づいて作成される対象プロファイル行列の2つのプロファイル行列を用意する手段と、

(b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出する手段と、

(c) 前記相関係数からなるスコア行列を作成する手段とを含むシステム。

【0026】

(2) (1) 記載のシステムにより作成されたスコア行列を用いることを特徴とするタンパク質立体構造の予測システム。

【0027】

(3) コンピュータを、タンパク質の立体構造を予測するためのタンパク質プロファイル行列間の類似性を評価するシステムとして機能させるためのプログラムであって、

前記プロファイル行列は、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントにおいて、各アミノ酸残基位置におけるアミノ酸種毎の出現確率を備えたプロファイルカラムの群から構成され、

前記類似性評価システムは、以下の手段：

(a) 立体構造を予測したいタンパク質を含む複数のタンパク質に基づいて作成される入力プロファイル行列と、立体構造が既知である複数のタンパク質に基づいて作成される対象プロファイル行列の2つのプロファイル行列を用意する手段と、

(b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出する手段と、

(c) 前記相関係数からなるスコア行列を作成する手段とを含むプログラム。

【0028】

(4) 上記(3) 記載のプログラムを記録したコンピュータ読み取り可能な記録媒体。

【0029】

(5) タンパク質プロファイル行列間の類似性を評価する方法であって、

前記プロファイル行列は、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントにおいて、各アミノ酸残基位置におけるアミノ酸種毎の出現確率を備えたプロファイルカラムの群から構成され、

前記類似性評価方法は、以下のステップ：

(a) 入力プロファイル行列と、対象プロファイル行列の2つのプロファイル行列を用意するステップと、

(b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出するステップと、

(c) 前記相関係数からなるスコア行列を作成するステップとを含む方法。

【0030】

(6) 前記対象プロファイル行列が、立体構造が既知である複数のタンパク質に基づいて作成されるプロファイル行列であり、前記入力プロファイル行列が、立体構造を予測したタンパク質を含む複数のタンパク質に基づいて作成されるプロファイル行列である上記

(5) 記載の類似性評価方法。

【0031】

(7) 上記(5)又は(6)で得られたスコア行列を用いることを特徴とするタンパク質立体構造の予測方法。

【発明の効果】

【0032】

本発明により、タンパク質プロファイル行列間の類似性を簡便かつ精度よく評価することができる。本発明により得られたスコア行列は、タンパク質立体構造を予測するのに好適に使用される。

【0033】

以下、本発明を詳細に説明する。

【発明を実施するための最良の形態】

【0034】

1. 類似度評価システム

図1は、本発明の一実施形態において使用されるハードウェア構成を示す図である。

【0035】

図1に示すように、本発明の類似性評価システムは、CPU101、ROM102、RAM103、入力部104、情報通信送信/受信部105、出力部106、ハードディスクドライブ(HDD)107及びCD-ROMドライブ108等を備える。

【0036】

CPU101は、情報記憶手段(例えば磁氣的及び/又は光学的記録媒体)に記憶されているプログラムに従って、類似性評価システム全体を制御する。そして、入力部104などから受け取った情報を出力部106に供給する。また、ネットワーク回線109を通じて受け取った情報に基づいて評価処理を実行することもできる。入力部104は、キーボードやマウス等であり、評価処理を実行する上で必要な条件又はデータを入力するときに操作される。ROM102は、本発明の類似性評価システムの動作に必要な処理を命令するプログラム等を格納する。RAM103は、類似性評価システムにおける処理を実行する上で必要なデータを一時的に格納する。

【0037】

送信/受信部105は、CPU101の命令に基づいて、ネットワーク回線109等との間で情報通信(データの送受信処理)を実行するものであり、例えばモデム、ルーター等が例示される。出力部106は、入力手段104から入力されたプロファイルデータ、その他各種条件等を、CPU101からの命令に基づいて情報表示処理する(例えば表示画面、プリンタ)。CD-ROMドライブ108は、CPU101の指示に基づいて、CD-ROMに格納されている類似性評価システムを機能させるためのプログラム又はデータ等を読み出し、例えばRAM103に格納する。CD-ROMの代わりに記録媒体として書き換え可能なCD-R、CD-RWを用いることもできる。その場合には、CD-ROMドライブ108の代わりにCD-R又はCD-RW用ドライブを設ける。また、上記媒体の他に、DVD、MOとそれらの媒体を用い、それに対応するドライブを備える構成としてもよい。

【0038】

コンピュータに本発明の類似性評価システムを機能させるためのプログラムは、例えば

C言語等で書くことができる。従って、このソフトウェアはWindows（登録商標）95/98/2000、Linux（登録商標）、UNIX（登録商標）等の各種オペレーティングシステムで作動させることが可能である。

【0039】

図2は、本発明のプロファイル行列間類似性評価システムを含む処理手順の一例を示すフローチャートである。

図2に示すように、本発明にかかる類似性評価システムでは、まず、比較したい2つのプロファイル行列（入力プロファイル行列と対象プロファイル行列）を用意し、続いてこれらの類似性を評価し、必要に応じて評価結果を出力する。以下、各処理について詳細に説明する。

【0040】

(a) プロファイル行列の用意 (S10)

プロファイル行列を用意するステップでは、比較したい2つのプロファイル行列が用意（抽出）される（S11、S12）。ここで、2つのプロファイル行列のうち、一方（対象プロファイル行列）は、立体構造が既知である複数のタンパク質に基づいて作成されたプロファイル行列（図2中、長さm）である。他方（入力プロファイル行列）は、立体構造を予測したいタンパク質（立体構造は未知であると既知であるとを問わない）を含む複数のタンパク質に基づいて作成されたプロファイル行列（図2中、長さn）であることが好ましい。

【0041】

プロファイル行列の作成方法としては、上述した従来知られている方法を採用することができ、特に制限はない。たとえば、ある一配列を入力値として、既存のプログラムであるPSI-BLASTを用いて、配列データベースに検索をかけてマルチプルアライメントを作成し、このマルチプルアライメントに基づいてプロファイル行列を作成してもよい。また、生物学的に相互に関連しあう複数のタンパク質のアミノ酸配列の一群を入力値として、既存のプログラムであるCLUSTALWを用いてマルチプルアライメントを作成し、当該マルチプルアライメントに基づいてプロファイル行列を作成してもよい。また、予め作成されたマルチプルアライメントを入力値とし、このマルチプルアライメントに基づいて作成してもよい。

【0042】

ここで、プロファイル行列は、ある代表アミノ酸配列の全配列に基づいて作成されていてもよく、また、代表配列中のモチーフ領域等、一部の領域に基づいて作成されていてもよい。また、マルチプルアライメントを作成する際に、経験的に導出されたギャップペナルティーを導入してもよい。

また、必要に応じて、プロファイル行列として、アミノ酸種の出現頻度を、アミノ酸種のランダム出現頻度で割った行列（PSSM: Gribskov, M., et al., (1987) Proc. Natl. Acad. Sci. USA, 84, 4355-4358）を用いてもよい。

【0043】

入力プロファイル行列は、たとえば、立体構造を予測したいタンパク質を代表アミノ酸配列として、この配列に基づいて作成することができる。また、対象プロファイル行列については、たとえば、SCOP (Murzin et al., J. Mol. Biol. 247(4):536-540 (1995)) やCATH (Orengo et al., Structure 5(8):1093-1108 (1997)) といったタンパク質構造分類データベースから取得したタンパク質のアミノ酸配列を代表配列とし、この配列に基づいて作成することができる。こうして得られた対象プロファイル行列は、代表配列ごとに予め作成しておき、対象プロファイル行列データベースとして保持しておくことが好ましい。

【0044】

(b) 相関係数の算出（プロファイル行列の比較評価） (S20)

続いて、プロファイル行列の類似性評価ステップでは、上記のステップで用意した入力プロファイル行列の各プロファイルカラムと、対象プロファイル行列の各プロファイルカ

ラムとの間の類似性を、各カラムペア毎に評価をする。

【0045】

図3は、各プロファイルカラムペア毎に類似性を評価し、スコア行列を作成するステップを模式的に示した図である。

本発明において、プロファイルカラム間の類似性は、プロファイルカラム間の相関係数を算出することによって行う。

【0046】

たとえば、入力プロファイル行列を $X = x_1 \ x_2 \ \cdots \ x_p \ \cdots \ x_n$ (ただし、 x_p はアミノ酸残基位置 p におけるプロファイルカラム) とし、対象プロファイル行列を $Y = y_1 \ y_2 \ \cdots \ y_q \ \cdots \ y_m$ (ただし、 y_q はアミノ酸残基位置 q におけるプロファイルカラム) としたときに、プロファイルカラム x_p および y_q 間の類似性スコア c_{pq} は、下記の式によって与えられる。

【0047】

【数2】

$$C_{pq} = \frac{\sum_a^j (x_{pa} - \overline{x_p})(y_{qa} - \overline{y_q})}{\sqrt{\sum_a^j (x_{pa} - \overline{x_p})^2 \sum_a^j (y_{qa} - \overline{y_q})^2}}$$

[式中、 x_{pa} = プロファイルカラム x_p の要素

y_{qa} = プロファイルカラム y_q の要素

$\overline{x_p}$ = プロファイルカラム x_p の平均値

$\overline{y_q}$ = プロファイルカラム y_q の平均値

j = プロファイルカラムの要素数 (通常 20) である。]

【0048】

本発明では、プロファイルカラム間の類似性をプロファイルカラム間の相関係数によって評価する。このため、プロファイルカラム間の相関の程度によって、類似性スコアが +1 から -1 の値をとることになる。たとえば、2つのプロファイルカラム中の要素間に相関がある場合、即ちアミノ酸置換パターンの傾向に類似性が有る場合には、相関係数は +1 に近い数値を取るようになる。また、2つのプロファイルカラムの各要素が互いにランダムな値を取っている場合、即ちアミノ酸置換パターンの傾向に相関が無い場合、相関係数は 0 になり、アミノ酸置換パターンの傾向が全く反対の場合、相関係数は -1 になり、アミノ酸置換パターンの傾向性の類似-非類似を非常に自然な形で表現する事が出来る。

【0049】

したがって、本発明では、アミノ酸残基の保存性が高い保存領域のような相関が高い領

域では、高い類似性スコアが得られるため、保存領域の類似性を精度よく評価することができる。

【0050】

また、本発明によれば、アミノ酸残基の保存性だけでなく、内積によって類似性を評価する従来の方法 (Rychlewski et alら) では不可能であった領域に関する類似性評価、たとえば、モチーフ内に存在する非保存位置や、タンパク質立体構造上露出していることが重要で極性のみが重大な意義を占める位置、あるいはその逆に埋没部分に位置し疎水性のみが保存されている位置といった、激しいアミノ酸置換があるもののその置換パターンに共通性があると考えられる領域についての類似性をより精度良く評価することが可能である。

【0051】

例えば、ある zinc fingerモチーフを有する2つのプロファイル行列を比較した場合を考えたとする。そのモチーフは

C-[DES]-x-C-x(3)-I

と表記される。これは、1, 4, 8番目の残基にそれぞれC, C, Iの残基が保存されており、2番目の残基では、D又はE又はSが出現し、3番目および、5, 6, 7番目の残基では保存残基が特に無いことが表されている。内積によって類似性を評価する従来の方法では、この場合、1, 2, 4, 8番目の残基位置では、高い数値を与えるが、その他の位置では低い数値しか与えない。したがって、内積によって類似性を評価する従来の方法は、モチーフの一部については類似性を評価しているものの、モチーフ全体の類似性については精度よく評価されていないということになる。

【0052】

しかしながら、本発明によれば、1, 2, 4, 8番目の残基位置に高い数値を与えるだけでなく、3, 5, 6, 7番目の残基位置においても、保存残基が特に無いという置換パターンの類似性を評価することが可能で、これら残基位置でも高い数値を与える。したがって、本発明によれば、モチーフ全体としてのパターン情報の全てを評価することが可能となる。

なお、本発明における類似性評価システムは、モチーフ領域に限られず、立体構造を予測したいタンパク質の配列全体に適用することができる。すなわち、ギャップペナルティを導入して得られたプロファイル行列間の類似性評価にも、好適に適用することができる。

【0053】

さらに、本発明によれば、スコア行列の各要素 (類似性スコア) の平均値および標準偏差がほぼ一定値をとるため、類似性スコアに対する煩雑な正規化処理を施す必要がないというメリットもある。

【0054】

(c) スコア行列の作成

プロファイルカラム間の相関係数 (類似性スコア) は、各プロファイルカラムの全部又は一部の組合せについて算出され、これに基づいてスコア行列が作成される。スコア行列は、類似性スコアが各プロファイルカラムの全組合せについて算出された場合は、入力プロファイル行列の長さを行とし、対象プロファイル行列の長さを列とする行列であり、類似性スコアが各プロファイルカラムの一部の組合せについて算出された場合は、その組合せの数に応じた行と列を持つ行列となる。

【0055】

図2の例では、類似性スコアは各プロファイルカラムの全組合せについて算出されており、入力プロファイル行列の長さがn、対象プロファイル行列の長さがmであることから、類似性スコアは $m \times n$ 個生成される (S22)。したがって、スコア行列はn行m列となる。スコア行列は、比較したいプロファイル行列の長さ、及び算出される類似性スコアの数に応じた行列を予め定義し (S21)、定義された行列の各カラムに、各プロファイルカラム間の相関係数を入力することにより作成することができる (S23)。

【0056】

本発明で得られたスコア行列によって、2つのプロファイル行列の最終スコア（行列間の類似性）を精度よく算出することができる。最終スコアは既知の手法により作成することができる。たとえば、図2の例では、比較されるプロファイル行列のそれぞれの代表アミノ酸配列と、本発明によって得られたこれらのプロファイル行列間のスコア行列を入力値として、ダイナミックプログラミングを用いて最適パスを算出する（S24）ことによって最終スコアを求めることができる（S25）。

【0057】

以上の操作を、対象プロファイル行列データベースに保持してある対象プロファイル行列のすべてに対して行うことが好ましい。

【0058】

2. タンパク質立体構造の予測（S30）

対象プロファイル行列ごとに得られた最終スコアは、タンパク質立体構造を予測するのに好適に使用される。たとえば、以下の既知の手順にしたがって処理をされる。

【0059】

(1) 入力値

まず、予測対象配列を含む入力プロファイル行列と、立体構造が既知である代表アミノ酸配列を含む対象プロファイル行列との最終スコア、および各代表配列の長さが入力される。このとき、対象プロファイル行列データベース中にN本の既知代表配列があれば、N個の最終スコアと配列長が入力されることになる。

【0060】

(2) 最終スコアの長さ依存性の補正

予測対象配列を含む入力プロファイル行列と、各既知代表配列を含む対象プロファイル行列との最終スコアは、代表配列長に依存した関係が認められる為、次のような統計処理を行う。まず、X軸に各代表配列の長さの自然対数をとった値、Y軸に予測対象配列を含む入力プロファイル行列と各既知代表配列を含むプロファイル行列との最終スコアをプロットし、異常なはずれ値を除いて回帰直線を引く。各長さ（即ちX軸でのある値）における平均値は回帰直線で表されるものとみなし、予測対象配列を含む入力プロファイル行列と各既知代表配列を含む対象プロファイル行列との最終スコアは、平均値からのずれで評価される。通常良く使用されるように、標準偏差を単位として、そのずれの度合いが測定される。

【0061】

(3) ソート

平均値からのずれが（高得点側に）大きいもの程類似性が有るとみなされる。それ故、平均値からのずれが（高得点側に）大きい順にソートされ、予測構造の候補とされる。

【0062】

(4) 予測構造としてのアライメントとスコア出力

上でソートされた順に予測構造の候補として出力される。結果全てを出力するのは無意味なため、予測精度を考慮し経験的に求められた閾値以上の平均値からのずれを有する結果のみを出力する。この時、予測精度の指標として、標準偏差を単位として計算される平均値からのずれの度合いが表示される。

【0063】

予測対象配列を含む入力プロファイル行列と、各既知代表配列を含む対象プロファイル行列とのアラインメントおよび最終スコアの結果は、ダイナミックプログラミングを用いて逐次計算された際のもので出力する。各既知代表配列は立体構造既知なので、このアラインメント出力が立体構造予測結果に相当する。

【0064】

3. コンピュータプログラム

本発明は、コンピュータを、タンパク質の立体構造を予測するためのタンパク質プロファイル行列間の類似性を評価するシステムとして機能させるためのプログラムをも提供する。本発明のコンピュータプログラムは、以下の手段：

- (a) 入力プロファイル行列と、対象プロファイル行列の2つのプロファイル行列を用意する手段と、
- (b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出する手段と、
- (c) 前記相関係数からなるスコア行列を作成する手段とを含むものである。

本発明のプログラムには、上記必須の手段以外に、汎用のプログラムとして通常備えられる汎用手段を含んでもよい。そのような手段としては、各種データの格納手段、情報の送受信手段、ディスプレイ、プリンター等の表示・出力手段等を挙げることができる。

【0065】

4. コンピュータ用記録媒体

本発明のプログラムは、コンピュータ読み取り可能な記録媒体又はコンピュータに接続しうる記憶手段に保存することができる。本発明のプログラムを含有するコンピュータ用記録媒体又は記憶手段も本発明に含まれる。記録媒体又は記憶手段としては、磁気的媒体（フレキシブルディスク、ハードディスクなど）、光学的媒体（CD、DVDなど）、磁気光学的媒体（MO、MD）などが挙げられる。

【実施例】

【0066】

以下、実施例により本発明をさらに具体的に説明する。但し、本発明はこれら実施例に限定されるものではない。

【0067】

実施例 1

(1) 対象プロファイル行列データベースの構築

構造分類データベース SCOP (URL: <http://scop.mrc-lmb.cam.ac.uk/scop/>) release 1.59 に基づく分類から、代表配列を取得した。その中から、単独ドメインを有し解像度 2.5 Å 以内の構造データを有するタンパク質のアミノ酸配列 948 本を選択した。948 本の代表配列各々に対して PSI-BLAST とアミノ酸配列データベース (NRDB: <ftp://ftp.ncbi.nlm.nih.gov>) より取得) を用いて対象プロファイル行列を構築し、対象プロファイル行列データベースを完成させた。

【0068】

ここで使用した「NRDB」には、現在知られているほぼ大部分のタンパク質アミノ酸配列が含まれている。PSI-BLAST を使うことで、この NRDB から各代表配列に生物学的に関連あると考えられる配列を自動的に収集し、さらにプロファイル行列も作成することが出来る。

【0069】

(2) 入力プロファイル行列の作成

本発明にかかるシステムによって正しい構造予測がなされているかどうかを調べるため、予測対象配列として構造が既に知られている配列、すなわち、対象プロファイル行列を作成する際に使用した上記 948 本の代表配列を使用した。入力プロファイル行列は、これらの予測対象配列を順次使用して、対象プロファイル行列の場合と同様の操作、すなわち、PSI-BLAST とアミノ酸配列データベース (NRDB) を用いて構築した。

【0070】

(3) 各プロファイル行列間の比較

続いて、上記で構築された予測対象配列（本実施例では 948 本の各代表配列）を含む入力プロファイル行列と、対象プロファイル行列データベース中の対象プロファイル行列との比較が順次なされた。この際、プロファイル行列間のスコア行列の各要素（類似性スコア）は、相関係数を用いて計算された。

こうして得られたプロファイル行列間のスコア行列を入力値として、ダイナミックプログラミングによってプロファイル行列間の最終スコアとアラインメントが出力された。

【0071】

各入力プロファイル行列に対して、以上の操作を対象プロファイル行列データベースに構築されたすべての対象プロファイル行列について行った。

【0072】

(4) 最終処理及び結果出力

評価の出力は、既に説明した方法に従って、948予測について各々結果出力を行った。すなわち、入力プロファイル行列と対象プロファイル行列との各最終スコアおよび各代表配列の長さを入力し、最終スコアの長さ依存性の補正を行った。続いて、平均値からのずれが(高得点側に)大きい順にソートし、ソートされた順に予測構造の候補として出力した。

こうして出力された予測構造の候補と、既にわかっている正しい予測構造とを比較することにより、予測結果の信頼度と感度を算出し、この結果を図4に示した。

【0073】

比較例1

実施例1で取得した948本の代表配列を用いて、配列類似性検索として一般的に用いられているPSI-BLASTを用いて構造予測を行った。すなわち、948本の代表配列各々に対してPSI-BLASTとアミノ酸配列データベース(NRDB:ftp://ftp.ncbi.nlm.nih.govより取得)を用いて構築したプロファイル行列を入力値とし、948本の代表配列に対して類似性検索を行い、予測構造の候補を出力した。

こうして出力された予測構造の候補と、既にわかっている正しい予測構造とを比較することにより、予測結果の信頼度と感度を算出し、この結果を図4に示した。

【0074】

比較例2

実施例1で取得した948本の代表配列を用いて、配列類似性検索として一般的に用いられているIMPALA(Schaffer, A. A., Wolf, Y. I., Ponting, C. P., Koonin, E. V., Aravind, L., and Altschul, S. F. (1999) Bioinformatics. 015:1000-1011)を用いて構造予測を行った。すなわち、948本の代表配列を入力値とし、948本の代表配列各々に対して予め作成し構築したプロファイル行列データベース(実施例1で構築した対象プロファイル行列データベースを使用した)に対して類似性検索を行い、予測構造の候補を出力した。

こうして出力された予測構造の候補と、既にわかっている正しい予測構造とを比較することにより、予測結果の信頼度と感度を算出し、この結果を図4に示した。

【0075】

図4から、比較例1および2の手法に比べて、信頼度0.98以降において、本発明にかかる実施例1が常に感度で勝っていることが示される。

【0076】

比較例3

プロファイル行列間のスコア行列の各要素(類似性スコア)を、内積法(Rychlewski et al. (2000), 9:p232-241)を用いて計算した以外は実施例1と同様の手法で予測構造の候補を出力した。

こうして出力された予測構造の候補と、既にわかっている正しい予測構造とを比較することにより、予測結果の信頼度と感度を算出し、この結果を図5に示した。

【0077】

実施例2

(1) 対象プロファイル行列データベースの構築

配列は、構造分類データベースSCOP(URL:http://scop.mrc-lmb.cam.ac.uk/scop/) release1.59に基づく分類から、お互いの同一残基率が40%未満であるドメイン単位の代表配列4381本を、SCOPの配列データベースであるASTRAL(http://astral.stanford.edu/)データベースから取得した。更に、タンパク質立体構造データベースPDB(URL:http://www.rcsb.org/pdb/)に登録されているが、SCOPに未登録であるものであって、ASTRALから取得

した上記4381本の配列と非類似のものを下記 (A) ~ (D) の要領で取得し、代表配列に加えた。このようにして選択されたアミノ酸配列各々に対して、下記 (A) ~ (D) の要領でPSI-BLASTとNRDBを用いて対象プロファイル行列を構築し、対象プロファイル行列データベースを完成させた。

【0078】

(A) 対象プロファイル行列データベースAの構築

2002年5月18日時点でのPDB中のアミノ酸配列をSCOPrelease1.59の分類に基づく代表配列に対してBLASTP(Altschul et al., Nucleic Acids Res. (1997) 25(17): 3389-3402: 非特許文献2)をかけ、期待値が0.00001以上のものを選んだ。さらにそれらを配列のクラスタリングを行うプログラムであるblastclustにかけ、互いの同一残基率が40%未満となるように配列248本を選択した。このようにして選択された配列と、SCOPrelease1.59の分類に基づく代表配列4381本との合計4629本の配列各々に対して、PSI-BLASTと2002年5月18日時点のNRDBを用いて対象プロファイル行列を構築し、対象プロファイル行列データベースAを完成させた。

【0079】

(B) 対象プロファイル行列データベースBの構築

2002年6月23日時点でのPDBと2002年5月18日時点でのPDB中のアミノ酸配列の差分を上記 (A) で作成した代表配列に対してBLASTPをかけ、期待値が0.00001以上のものを選んだ。さらにそれらをblastclustにかけ、互いの同一残基率が40%未満となるように配列49本を選択した。このようにして選択された配列と、上記 (A) で作成した代表配列との合計4678本の配列各々に対して、PSI-BLASTと2002年6月17日時点のNRDBを用いて対象プロファイル行列を構築し、対象プロファイル行列データベースBを完成させた。

【0080】

(C) 対象プロファイル行列データベースCの構築

2002年7月14日時点でのPDBと2002年6月23日時点でのPDB中のアミノ酸配列の差分を上記 (B) で作成した代表配列に対してBLASTPをかけ、期待値が0.00001以上のものを選んだ。さらにそれらをblastclustにかけ、互いの同一残基率が40%未満となるように配列23本を選択した。このようにして選択された配列と、上記 (B) で作成した代表配列との合計4701本の配列各々に対して、PSI-BLASTと2002年7月9日時点のNRDBを用いて対象プロファイル行列を構築し、対象プロファイル行列データベースCを完成させた。

【0081】

(D) 対象プロファイル行列データベースDの構築

上記 (C) で作成した代表配列の合計4701本の配列各々に対して、PSI-BLASTと2002年8月6日時点のNRDBを用いて対象プロファイル行列を構築し、対象プロファイル行列データベースDを完成させた。

【0082】

(2) 入力プロファイル行列の作成

配列は、隔年で行われる世界的規模で行われる構造予測コンテストの2002年度大会であるCASP5/CAFASP3(URL:<http://predictioncenter.llnl.gov/casp5/>)において、構造認識部門(通常の配列解析手法では立体構造既知であるタンパク質と明白な配列類似性を有さないが、その構造が(実際に解かれてみると)既知立体構造との構造類似性を有する、即ち類似性検索が困難なタンパク質に関する予測する部門)において出題された配列、すなわち、現在通常の配列解析手法(例えば、PSI-BLASTなど)では、立体構造既知であるタンパク質と明白な配列類似性を有さないタンパク質であり、かつ、その構造が(実際に解かれてみると)既知立体構造との構造類似性が明らかになったアミノ酸配列を用いた。具体的には、URL:<http://www.cs.bgu.ac.il/~dfischer/CAFASP3/targets.html>において、下記のターゲット番号が付されたアミノ酸配列22本を用いた。

【0083】

T0130、T0132、T0134、T0135、T0136、T0138、T0146、T0147、T0148、T0156、T0157、T0159、T0162、T0168、T0170、T0172、T0173、T0174、T0186、T0187、T0191、T0193

【0084】

これら22本の配列各々に対して、PSI-BLASTとNRDBを用いて入力プロファイル行列を構築し、入力プロファイル行列データベースを完成させた。

なお、NRDBとしては、2002年5月18日時点、2002年6月17日時点、2002年7月9日時点、及び2002年8月6日時点のものの計4種類を使用し、得られた入力プロファイル行列データベースを、それぞれ、「入力プロファイル行列データベースA」、「入力プロファイル行列データベースB」、「入力プロファイル行列データベースC」、及び「入力プロファイル行列データベースD」とした。

【0085】

(3) 各プロファイル行列間の比較

続いて、上記で構築された予測対象配列を含む入力プロファイル行列データベースAの入力プロファイル行列と、対象プロファイル行列データベースA中の対象プロファイル行列との比較を、実施例1の「(3)各プロファイル行列間の比較」と同様の手順で行った（比較A）。

同様の操作を、入力プロファイル行列データベースBと対象プロファイル行列データベースBに対して、入力プロファイル行列データベースCと対象プロファイル行列データベースCに対して、及び、入力プロファイル行列データベースDと対象プロファイル行列データベースDに対して、それぞれ行った（比較B、C、D）。

【0086】

(4) 最終処理及び結果出力

評価の出力は、既に説明した方法に従って22予測について各々結果出力を行った。即ち、各データベースの組合せ（比較A～D）においてそれぞれ得られた、入力プロファイル行列と対象プロファイル行列との各最終スコアおよび、各代表配列の長さを入力し、最終スコアの長さ依存性を補正した。続いて平均値からのずれが、（高得点側に）大きい順にソートし、ソートされた順に上位10個までを予測構造の候補として22本の配列各々に対して出力した（出力A～D）。

こうして出力された予測構造の候補と、コンテストの予測構造投稿期間の後に公開された実験により解かれた立体構造とを比較することで、予測結果の正確さが測定された。予測構造評価方法の一つは、予測構造と正解構造の重ね合わせを行い、対応残基が3Åより短い距離にある残基数を出力A～Dについて積算すること（sum値）により行われた。22のタンパク質を構造ドメイン単位（全部で34ドメイン）で眺めた結果によれば、構造予測コンテストCASP5/CAFASP3における上記構造認識部門において22本の配列各々に対して上位1個の予測を考慮した時、本手法のsum値は「577」であり、これは、配列情報を用いた他のいかなる手法よりも優れているものであった。

また、ある閾値を設定してある入力（予測対象）配列に対する予測の成否を観測した場合でも、22本の配列各々に対して上位1個の予測を考慮した時本手法は、予測が成功したと判断される個数を出力A～Dについて積算したもの（correct値）において、「9」と高く、配列情報を用いた他のいかなる手法よりも優れていることが示された。

【図面の簡単な説明】

【0087】

【図1】本発明の一実施形態において使用されるハードウェア構成を示す図である。

【図2】本発明のプロファイル行列間類似性評価システムを含む処理手順の一例を示すフローチャートである。

【図3】本発明のプロファイル行列間類似性評価システムにおいて、各プロファイルカラムペア毎に類似性を評価し、スコア行列を作成するステップを示す図である。

【図4】実施例1、比較例1及び比較例2において出力された予測結果の信頼度と感度とをプロットした図である。

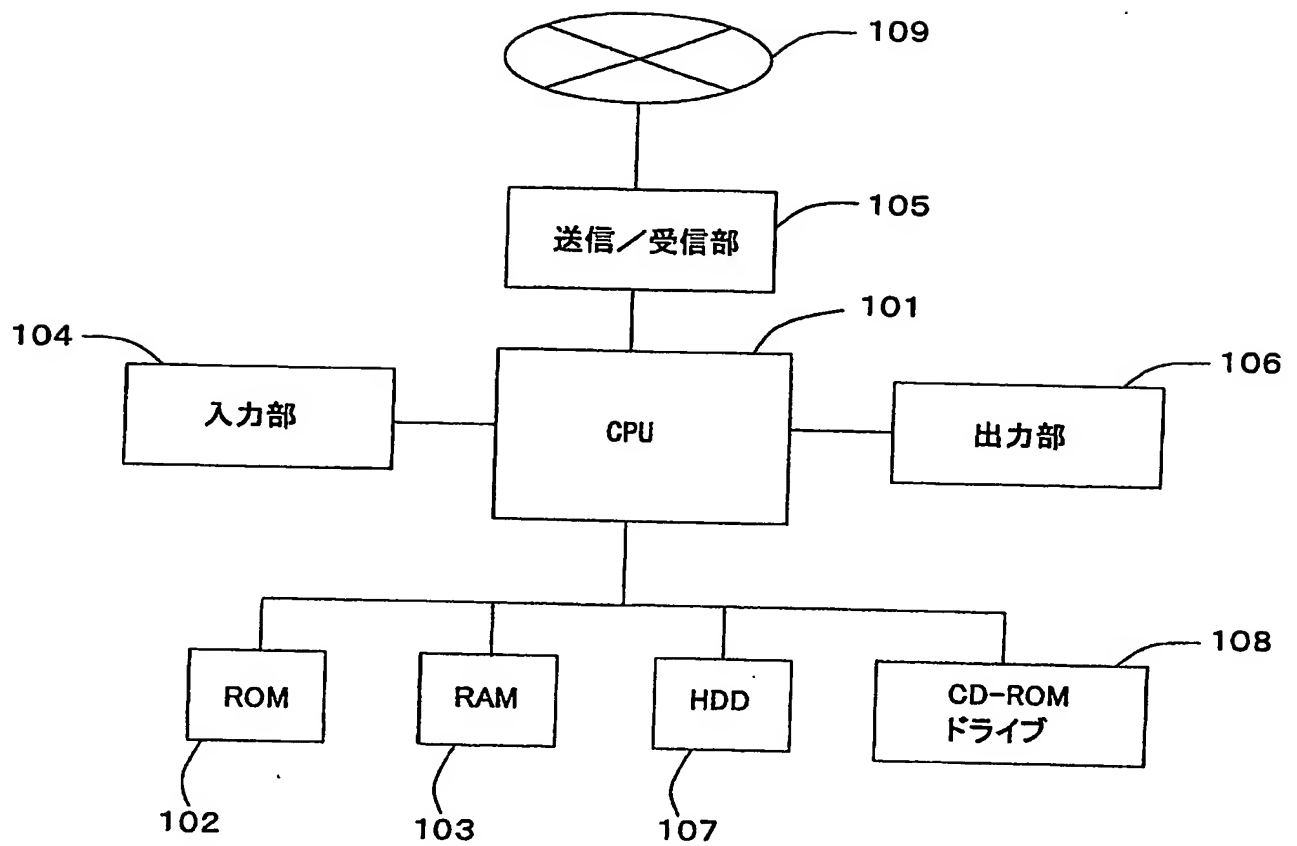
【図5】実施例1及び比較例3において出力された予測結果の信頼度と感度とをプロットした図である。

【符号の説明】

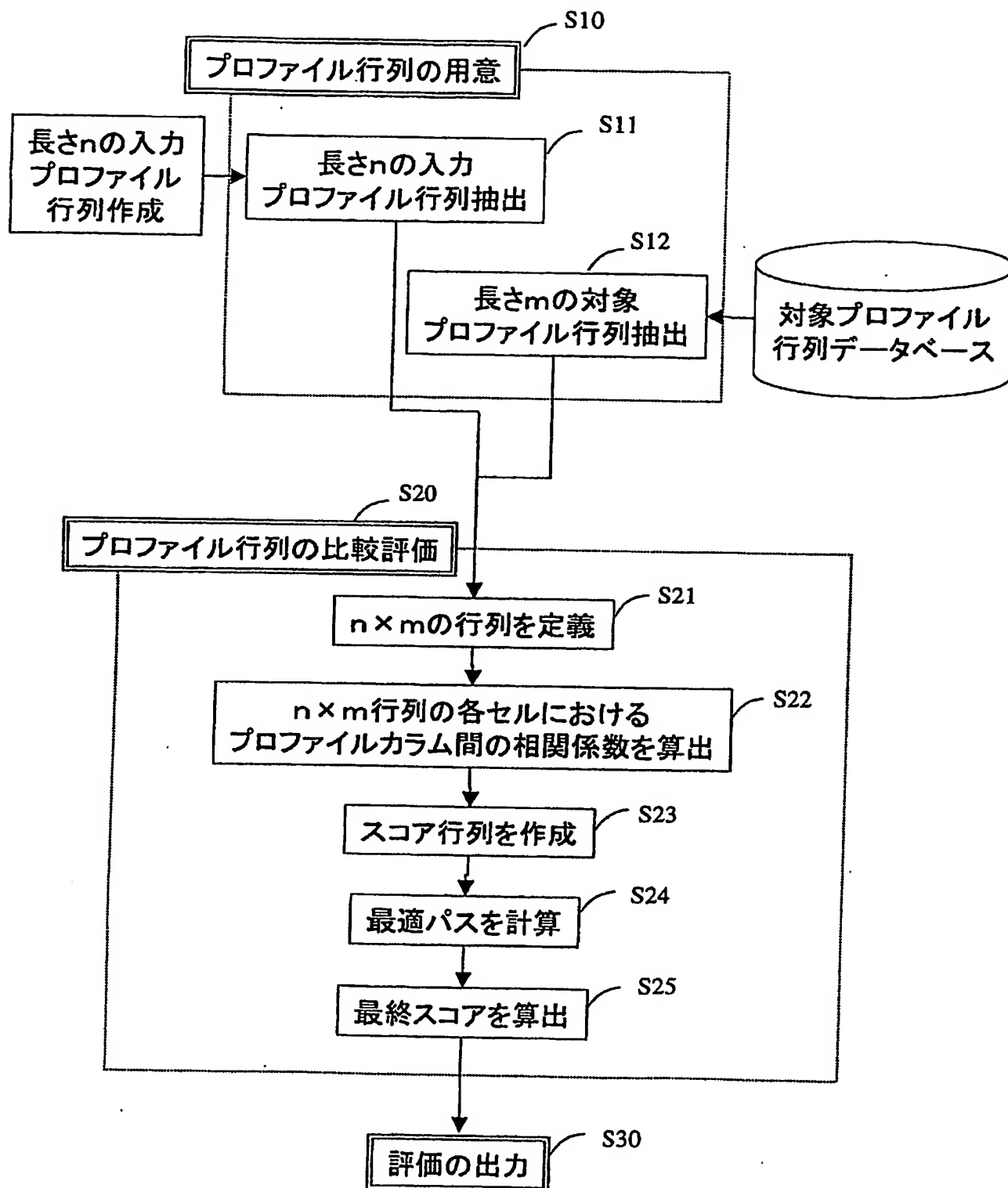
【 0 0 8 8 】

101：CPU、
102：ROM、 103：RAM、 104：入力部、105：送信/受信部、
106：出力部、 107：HDD、 108：CD-ROMドライブ、109：ネットワーク回線

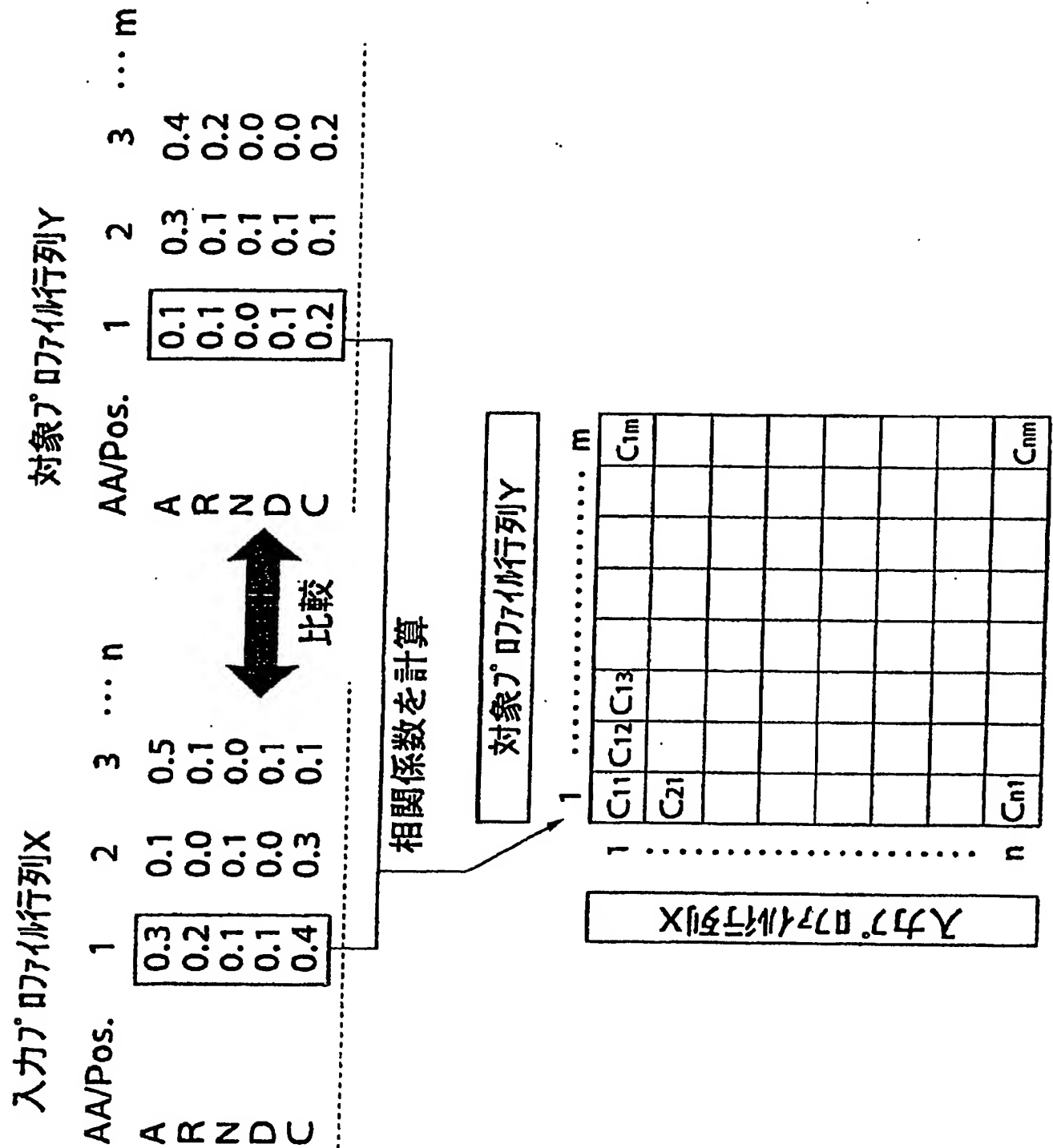
【書類名】 図面
【図 1】



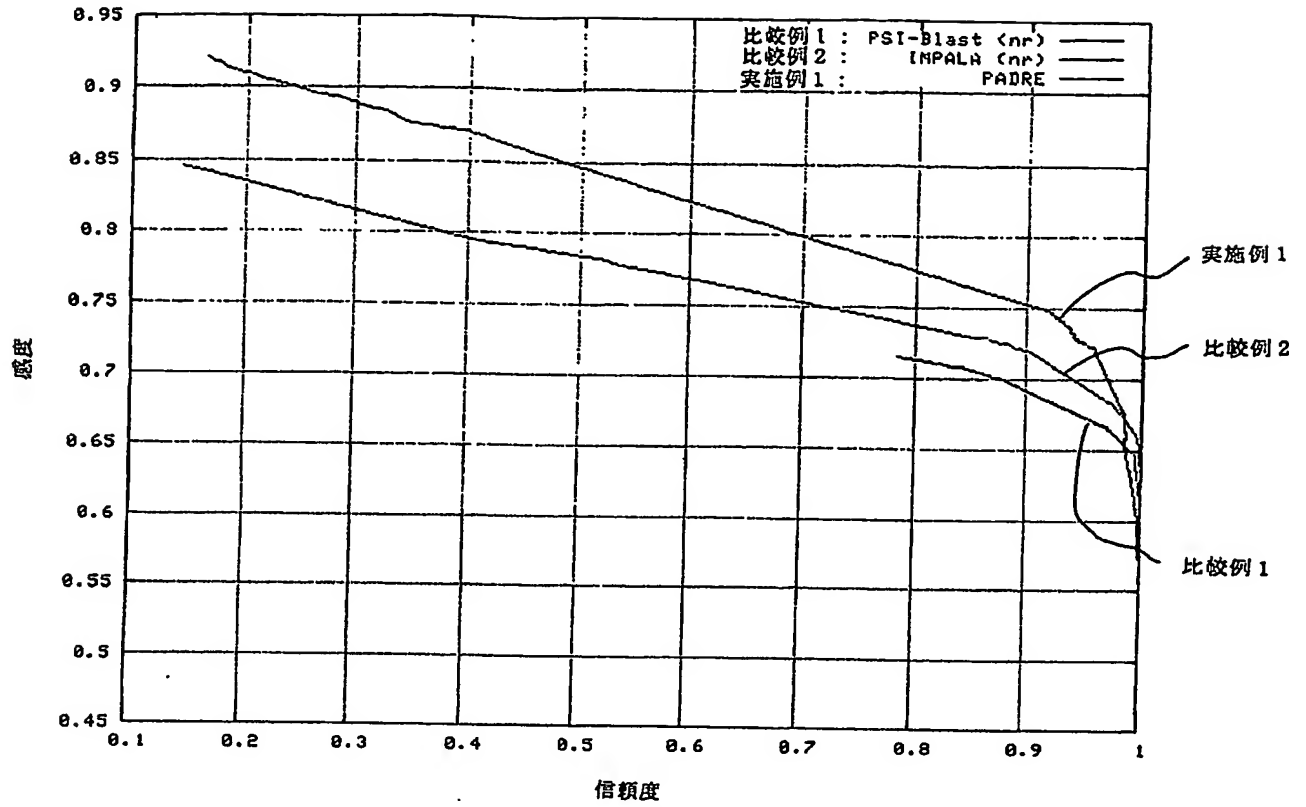
【図2】



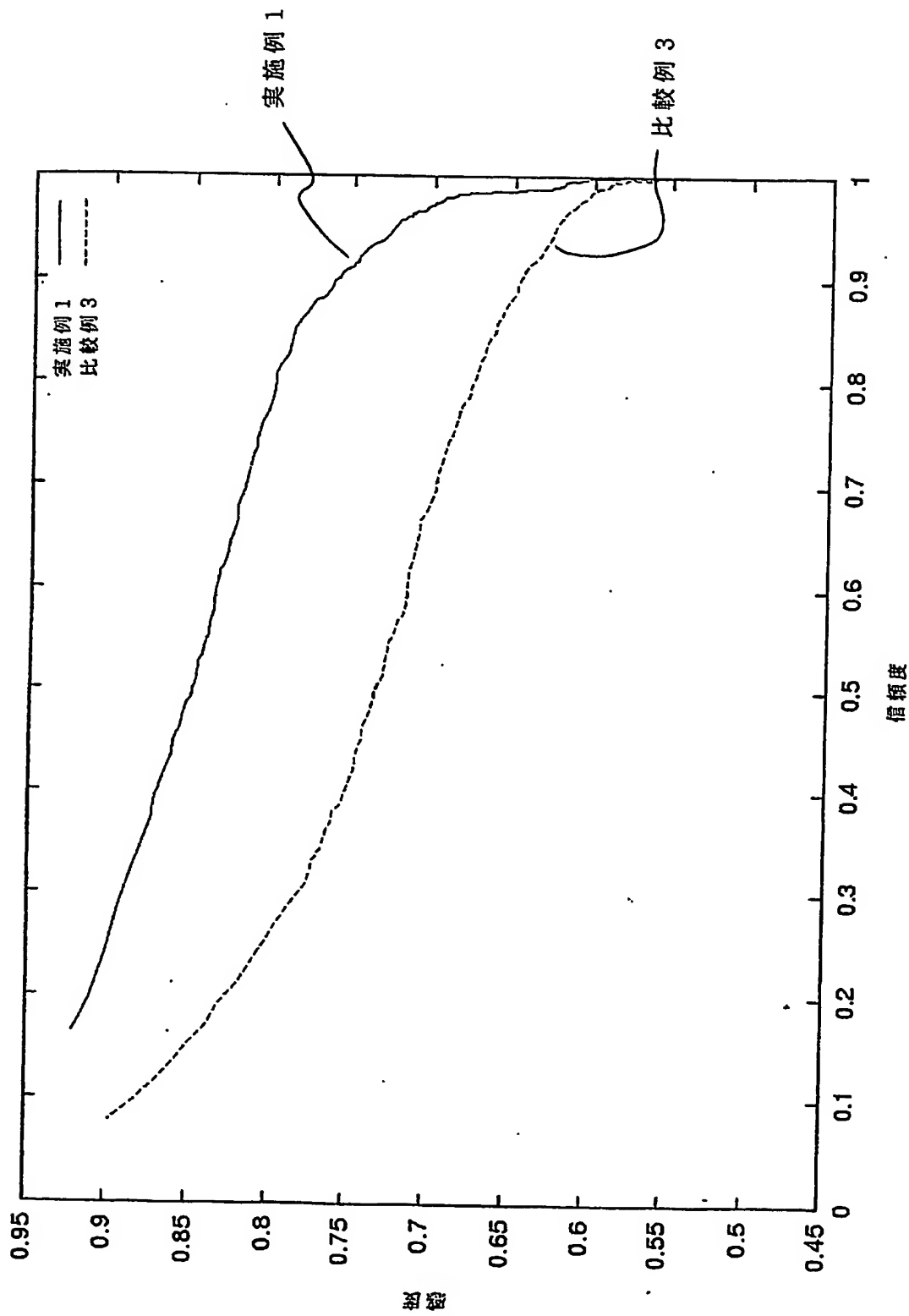
【図 3】



【図 4】



【図5】



【書類名】 要約書

【要約】

【課題】 タンパク質の立体構造予測に好適に使用できる、タンパク質プロファイル行列間の類似性評価システムの提供。

【解決手段】 タンパク質プロファイル行列間の類似性を評価するシステムであつて、プロファイル行列は、関連する複数のタンパク質のアミノ酸配列を多重並置させたマルチプルアライメントにおいて、各アミノ酸残基位置におけるアミノ酸種毎の出現確率を備えたプロファイルカラムの群から構成され、(a) 入力プロファイル行列と、対象プロファイル行列の2つのプロファイル行列を用意する手段と、(b) 前記入力プロファイル行列の各プロファイルカラムと、前記対象プロファイル行列の各プロファイルカラムとの間の相関係数を、各プロファイルカラムの全部又は一部の組合せについて算出する手段と、(c) 前記相関係数からなるスコア行列を作成する手段とを含むシステムにより、上記課題を解決する。

【選択図】 図2

認定・付加情報

特許出願の番号	特願 2003-406776
受付番号	50302005591
書類名	特許願
担当官	第一担当上席 0090
作成日	平成 15 年 12 月 10 日

<認定情報・付加情報>

【提出日】	平成15年12月 5日
【特許出願人】	申請人
【識別番号】	301021533
【住所又は居所】	東京都千代田区霞が関 1-3-1
【氏名又は名称】	独立行政法人産業技術総合研究所

特願 2003-406776

出願人履歴情報

識別番号

[301021533]

1. 変更年月日

2001年 4月 2日

[変更理由]

新規登録

住 所

東京都千代田区霞が関1-3-1

氏 名

独立行政法人産業技術総合研究所